# Coursera and Campuswire Query Link

CS410 Project

Xiaofan Liu / Boyu Pang / Brittany West

December 2021

# Introduction

# Project Process

- Collect data from Coursera
- Cleaning the dataset
- Build the data corpus
- Collect sample queries from Campuswire
- Generating query result judgement manually
- Ranking function implementation
- Testing, fine-tuning and improving ranking functions
- Output the result
- Documentation and presentation

# Libraries and Files

# Packages Needed

- Coursera-dl
- Metapy
- Html2text
- Numpy
- Pandas
- Regular Expression

# Files

- Original txt and html files
- DataCleaning.py
- Config.toml
- Documents.txt
- CampusWireHeading_Queries.txt
- Qrels.txt
- ModelEvaluation.py
- Search.py

# Code Implementation

# Data Cleaning

```
1    This lecture is the first one about the text clustering. In this lecture, we are going to talk
2    This lecture is about the Evaluation of Text Categorization. So we've talked about many differ
3    This lecture is about the Latent Aspect Rating Analysis for Opinion Mining and Sentiment Analy
4    In this lecture we give an overview of Text Mining and Analytics. First, let's define the term
5    . This lecture is about the syntagmatic relation discovery, and entropy. In this lecture, we're
6    This lecture is about the mixture model estimation. In this lecture, we're going to continue dis
7    Hello welcome to CS410 DSO Text Information Systems. This is an online course offered by Univers
8    ### **Introduction** The course project is to give the students hands-on experience on develop
9    ## Exam Instructions    * A password quiz precedes and unlocks the proctored exam. The proctor
10   >> This lecture is about Natural Language of Content Analysis. As you see from this picture, t
11   In this lecture, we are going to talk about how to improve the instantiation of the vector spa
12   This lecture is about Evaluation of Text Retrieval Systems In the previous lectures, we have t
13   This lecture is about the Probabilistic Retrieval Model. In this lecture, we're going to conti
14   This lecture is about the feedback in text retrieval. So in this lecture, we will continue wit
15   This lecture is about the Learning to Rank. In this lecture, we are going to continue talking
16   No office hour this week due to Fall Break
17   In this course, there are two timed exams proctored via [ ProctorU ](https://www.proctoru.com/pc
18   ## **CS410 Technology Review (4-credit students only)** **CS410 Technology Review** The Techn
19   # Week 1 Overview  The first six weeks of the course are based on the content of the Text Retri
20   # Week 10 Overview  During this week's lessons, you will learn text clustering, including the b
21   # Week 11 Overview  During this week's lessons, you will continue learning about various methoc
22   # Week 2 Overview  During this week's lessons, you will learn how the vector space model works
23   # Week 3 Overview  During this week's lessons, you will learn how to evaluate an information re
24   # Week 4 Overview  During this week's lessons, you will learn probabilistic retrieval models ar
25   # Week 5 Overview  During this week's lessons, you will learn feedback techniques in informatic
26   # Week 6 Overview  During this week's lessons, you will learn how machine learning can be used
27   # Week 7 Overview  From Week 7 to Week 12, the lectures are based on the Text Mining and Analyt
```

```python
#Libraries for Data Collection and Cleaning
import os
import html2text
#Preprocess Data
os.remove("documents.txt")
c = open("documents.txt",'a+',encoding = "utf-8")
for filename in os.listdir('Data'):
    if filename.endswith('.txt'):
        with open(os.path.join('Data', filename)) as f:
            content = f.read()
            content = content.replace('\n',' ')
            content = content.replace('[NOISE]',' ')
            content = content.replace('[MUSIC]',' ')
            content = content.replace('[SOUND]',' ')
            content = content.replace('\u2011',' ')
            c.write(content + '\n')

    if filename.endswith('.html'):
        with open(os.path.join('Data', filename), encoding='utf8') as d:
            content = d.read()
            h = html2text.HTML2Text()
            content = h.handle(content)
            content = content.replace('\n',' ')
            content = content.replace('\u2011',' ')
            string = str(content)
            c.write(string + '\n')
c.close()
```

# Implementation of Ranking Functions

```
58   cfg = "config.toml"
59   idx = metapy.index.make_inverted_index(cfg)
60
61   # testing code - not used in final program
62   with open(cfg, 'r') as fin:
63       cfg_d = pytoml.load(fin)
64
65   query_cfg = cfg_d['query-runner']
66   if query_cfg is None:
67       print("query-runner table needed in {}".format(cfg))
68       sys.exit(1)
69
70   top_k = 10
71   query_path = query_cfg.get('query-path', 'CampusWireHeading_Queries.txt')
72   query_start = query_cfg.get('query-id-start', 0)
73
74   # get the query from user input
75   query_text = input("please enter a search query: ")
76   query = metapy.index.Document()
77   query.content(query_text.lower())
```
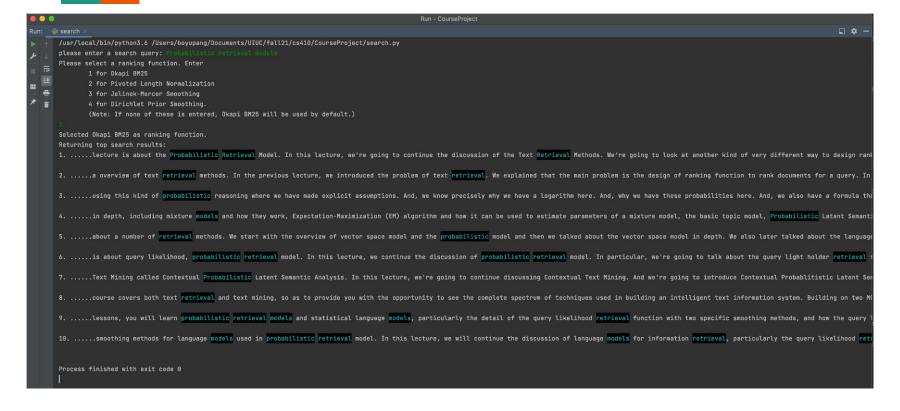
```
82    # get ranking function from user input
83    ranking_function = input("""Please select a ranking function. Enter
84        1 for Okapi BM25
85        2 for Pivoted Length Normalization
86        3 for Jelinek-Mercer Smoothing
87        4 for Dirichlet Prior Smoothing.
88        (Note: If none of these is entered, Okapi BM25 will be used by default.)
89  """)
90
91
92    if ranking_function == "2":
93        print("Selected Pivoted Length Normalization as ranking function.")
94        ranker = metapy.index.PivotedLength(s=0.2)
95    elif ranking_function == "3":
96        print("Selected Jelinek-Mercer Smoothing as ranking function.")
97        ranker = metapy.index.JelinekMercer(.85)
98    elif ranking_function == "4":
99        print("Selected Dirichlet Prior Smoothing as ranking function.")
100       ranker = metapy.index.DirichletPrior(2000)
101   else:
102       print("Selected Okapi BM25 as ranking function.")
103       ranker = metapy.index.OkapiBM25(k1=1.2, b=.75, k3=500)
104
105   top_docs = ranker.score(idx, query, num_results=10)
```

[(54, 1.2830500602722168), (40, 1.2738627195358276), (29, 1.2669488191604614), (63, 1.2639679908752441), (23, 1.2637956142425537), (98, 1.2576667070388794), (5, 1.2554749250411987), (12, 1.254915714263916), (31, 1.2459347248077393), (34, 1.2417280673980713)]

# Testing and Fine-tuning Ranking Functions

| | Ranker | k1 | b | k3 | s | lambda | mu | MAP | NDCG |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Okapi | 1.2 | 0.75 | 500 | NA | NA | NA | 0.81 | 0.85 |
| 1 | Okapi | 1.3 | 0.8 | 500 | NA | NA | NA | 0.79 | 0.84 |
| 2 | Okapi | 1.1 | 0.5 | 500 | NA | NA | NA | 0.77 | 0.84 |
| 3 | Pivoted Length | NA | NA | NA | 0.2 | NA | NA | 0.79 | 0.83 |
| 4 | Pivoted Length | NA | NA | NA | 0.01 | NA | NA | 0.79 | 0.83 |
| 5 | Pivoted Length | NA | NA | NA | 0.5 | NA | NA | 0.66 | 0.73 |
| 6 | JelineK-Mercer | NA | NA | NA | NA | .7 | NA | 0.66 | 0.79 |
| 7 | JelineK-Mercer | NA | NA | NA | NA | .5 | NA | 0.70 | 0.81 |
| 8 | JelineK-Mercer | NA | NA | NA | NA | .85 | NA | 0.70 | 0.81 |
| 9 | Dirichlet Prior | NA | NA | NA | NA | NA | 2000 | 0.66 | 0.78 |
| 10 | Dirichlet Prior | NA | NA | NA | NA | NA | 1500 | 0.66 | 0.78 |
| 11 | Dirichlet Prior | NA | NA | NA | NA | NA | 3500 | 0.66 | 0.78 |

# Output the Result

# Live Test